

Automatic Hidden-Web Table Interpretation, Conceptualization, and Semantic Annotation

Cui Tao*and David W. Embley*
Department of Computer Science
Brigham Young University, Provo, Utah 84602, U.S.A.

Abstract

The longstanding problem of automatic table interpretation still illudes us. Its solution would not only be an aid to table processing applications such as large volume table conversion, but would also be an aid in solving related problems such as information extraction, semantic annotation, and semi-structured data management. In this paper, we offer a solution for the common special case in which so-called sibling pages are available. The sibling pages we consider are pages on the hidden web, commonly generated from underlying databases. Our system compares them to identify and connect nonvarying components (category labels) and varying components (data values). We tested our solution using more than 2,000 tables in source pages from three different domains — car advertisements, molecular biology, and geopolitical information. Experimental results show that the system can successfully identify sibling tables, generate structure patterns, interpret tables using the generated patterns, and automatically adjust the structure patterns as it processes a sequence of hidden-web pages. For these activities, the system was able to achieve an overall F-measure of 94.5%. Further, given that we can automatically interpret tables, we next show that this leads immediately to a conceptualization of the data in these interpreted tables and thus also to a way to semantically annotate these interpreted tables with respect to the ontological conceptualization. Labels in nested table structures yield ontological concepts and interrelationships among these concepts, and associated data values become annotated information. We further show that semantically annotated data leads immediately to queryable data. Thus, the entire process, which is fully automatic, transform facts embedded within tables into facts accessible by standard query engines.

1 Introduction

The World Wide Web serves as a powerful resource for every community. Much of this online information, indeed, the vast majority, is stored in databases on the so-called hidden web.¹ Hidden-web information is usually only accessible to users through search forms and is typically presented to them in tables. Automatically understanding these hidden-web pages and making their facts externally accessible is a challenging task. In this paper, we introduce a domain independent, web-site independent, unsupervised way to automatically interpret tables from hidden-web pages.

*Supported in part by the National Science Foundation under Grant #0414644.

¹There are more than 500 billion hidden-web pages. The surface web, which is indexed by common search engines only constitutes less than 1% of the World Wide Web. The hidden web is several orders of magnitude larger than the surface web [28].

Once interpreted, we can automatically annotate the information in these pages and make it available to standard query engines.

Tables present information in a simplified and compact way in rows and columns. Data in one row/column usually belongs to the same category or provides values for the same concept. The labels of a row/column describe this category or concept.

Although a table with a simple row and column structure is common, tables can be much more complex. Figure 1 shows an example. Tables may be nested or conjoined as are the tables in Figure 1. Labels may span across several cells to give a general description as does *Identification* and *Location* in Figure 1. Although labels commonly appear on the top or left, table designers occasionally place labels on the right side of a table. In long tables, labels sometimes appear at the end of a table or in the middle of a table, every few rows, in order to help a reader find the correspondence between labels and data. Sometimes tables are rearranged to fit the space available. Label-value pairs may appear in multiple columns across a page or in multiple rows placed below one another down a page. These complexities make automatic table interpretation challenging.

To interpret a table is to properly associate table category labels with table data values. Using Figure 1 as an example, we see that *Identification*, *Location*, and *Function* are labels for the large rectangular table. Inside the right cell of the first row is another table with headers *IDs*, *NCBI KOGs*, *Species*, etc. Nested inside of this cell are two tables with labels *CGC name*, *Sequence name*, *Other name(s)*, *WB Gene ID*, *Version*, and *Gene Model*, *Status*, *Nucleotides (coding/transcript)*, *Protein*, and *Amino Acids*. Most of the rest of the text in the large rectangular table comprises the data values. If we look more closely, however, we may conclude that some category labels are interleaved in the text. For example, *via person* appears to be a label under *CGC name*, as does *Entrez Genes* and *Ace View* beside *NCBI*.

Once category labels and data values are found, we want to properly associate them. For example, the associated label for the value *F18H3.5* should be the sequences of labels *Identification*, *IDs*, and *Sequence name*. Given the source table in Figure 1, we match category labels with values as Figure 2 shows. We associate one or more sequences of labels with each data value in a table. Borrowing notation from Wang [45], the left hand side of the arrow is a sequence of one or more table labels, and the right hand side of the arrow is a data value. For the first two label-value pairs in Figure 2, there is only one label sequence. The third, however, has two: *Identification.Gene model(s).Amino Acids* and *2*. Each label sequence represents a dimension. In general, a table may have one, two, three, or more dimensions. If a table has multiple records (usually multiple rows) and if the records do not have labels, we add record numbers. The table under *Identification.Gene model(s)*, for example, has two records (two rows), but no row labels. We therefore label the first record *1* and the second record *2*. Thus, the label-value association becomes (*Identification.Gene*

Gene Summary for *cdk-4*

Specify a gene using a gene name ([unc-26](#)), a predicted gene id ([R13A5.9](#)), or a protein ID ([CE02711](#)):

[\[identification\]](#) [\[location\]](#) [\[function\]](#) [\[gene ontology\]](#) [\[reactome knowledgebase\]](#) [\[alleles\]](#) [\[similarities\]](#) [\[reagents\]](#) [\[bibliography\]](#)

Identification	IDs:	<table border="1"> <thead> <tr> <th>CGC name</th> <th>Sequence name</th> <th>Other name(s)</th> <th>WB Gene ID</th> <th>Version</th> </tr> </thead> <tbody> <tr> <td>cdk-4 - (Cyclin-Dependent Kinase family) (via person: Michael Krause)</td> <td>F18H3.5</td> <td>NM_077855 (inferred automatically) XO136 (inferred automatically)</td> <td>WBGene00000406</td> <td>1</td> </tr> </tbody> </table>	CGC name	Sequence name	Other name(s)	WB Gene ID	Version	cdk-4 - (Cyclin-Dependent Kinase family) (via person: Michael Krause)	F18H3.5	NM_077855 (inferred automatically) XO136 (inferred automatically)	WBGene00000406	1				
	CGC name	Sequence name	Other name(s)	WB Gene ID	Version											
	cdk-4 - (Cyclin-Dependent Kinase family) (via person: Michael Krause)	F18H3.5	NM_077855 (inferred automatically) XO136 (inferred automatically)	WBGene00000406	1											
	NCBI KOGs⁺:	Protein kinase PCTAIRE and related kinases [KOG0594]														
	Species:	<i>Caenorhabditis elegans</i>														
Other sequence(s):	AF083878 (Caenorhabditis elegans cyclin-dependent kinase CDK-4 (cdk-4) mRNA, complete cds.)															
NCBI Gene model(s):	[Entrez Genes: 15718266] [AceView: XO136]															
	<table border="1"> <thead> <tr> <th>Gene Model</th> <th>Status</th> <th>Nucleotides (coding/transcript)</th> <th>Protein</th> <th>Amino Acids</th> </tr> </thead> <tbody> <tr> <td>F18H3.5a 1, 2</td> <td>confirmed by cDNA(s)</td> <td>1029/3051 bp</td> <td>WP:CE18608</td> <td>342 aa</td> </tr> <tr> <td>F18H3.5b 1, 2, 3</td> <td>partially confirmed by cDNA(s)</td> <td>1221/1704 bp</td> <td>WP:CE28918</td> <td>406 aa</td> </tr> </tbody> </table>	Gene Model	Status	Nucleotides (coding/transcript)	Protein	Amino Acids	F18H3.5a 1, 2	confirmed by cDNA(s)	1029/3051 bp	WP:CE18608	342 aa	F18H3.5b 1, 2, 3	partially confirmed by cDNA(s)	1221/1704 bp	WP:CE28918	406 aa
Gene Model	Status	Nucleotides (coding/transcript)	Protein	Amino Acids												
F18H3.5a 1, 2	confirmed by cDNA(s)	1029/3051 bp	WP:CE18608	342 aa												
F18H3.5b 1, 2, 3	partially confirmed by cDNA(s)	1221/1704 bp	WP:CE28918	406 aa												
Putative ortholog(s):	<i>Caenorhabditis briggsae</i> : CBG07433 [syntenic alignment] (Stein LD et al. best reciprocal blastp match-seg-off)															
Location	Genetic Position: X:12.69 +/- 0.000 cM [mapping data] Genomic Position: X:13518823..13515773 bp															
Function	Mutant Phenotype: [Krause MW] cdk-4 is a cyclin dependent kinase related to cdk-4 and cdk-6 from other organisms. Homozygous cdk-4(gv3) animals usually arrest in L2 due to no, or limited, proliferation of the post-embryonic blast cells. About 3% of animals make it to a late stage of development. Definitions of abbreviations used in the text. RNAi Phenotype(s): Lvl Pvl Unc [For details see: Park M 06 Oct 1999]															

Figure 1: A Sample Table from WormBase [1].

$model(s).Amino\ Acids, 2) \rightarrow 406\ aa$ where $Identification.Gene\ model(s).Amino\ Acids$ is the label for the first dimension, and 2 is the row label for the second dimension.

Although automatic table interpretation can be complex, if we have another page, such as the one in Figure 3, that has essentially the same structure, the system might be able to obtain enough information about the structure to make automatic interpretation possible. We call pages that are from the same web site and have similar structures *sibling pages*.² The two pages in Figures 1 and 3 are sibling pages. They have the same basic structure, with the same top banners that appear in all the pages from this web site, with the same table title (*Gene Summary for some*

²Hidden-web pages are usually generated dynamically from a pre-defined templates in response to submitted queries, there fore they are usually sibling pages

(Identification.IDs.CGC name) →
cdk-4-(Cyclin-Dependent Kinase family)
(via person: Michael Krause);
(Identification.IDs.Sequence name) → *F18H3.5;*
 ...
(Identification.Gene model(s).Amino Acids, 2) → *406 aa;*
 ...

Figure 2: Interpretation for the Tables in Figure 1 (Partial).

particular gene), and a table that contains information about the gene. Corresponding tables in sibling pages are called *sibling tables*. If we compare the two large tables in the main part of the sibling pages, we can see that the first columns of each table are exactly the same. If we look at the cells under the *Identification* label in the two tables, both contain another table with two columns. In both cases, the first column contains identical labels *IDs*, *NCBI KOGs*, ..., *Putative orthology(s)*. Further, the tables under *Identification.IDs* also have identical header rows. The data rows, however, vary considerably. General speaking, we can look for commonalities to find labels and look for variations to find data values.

Given that we can find most of the label and data cells in this way, our next task is to infer the general structure pattern of the web site and of the individual tables embedded within pages of the web site. With respect to identified labels, we look below or to the right for value associations; we may also need to look above or to the left. In Figure 1, the values for *Identification.Gene Model(s).Gene Model* are below, and the values for *Identification.Speices* are to the right.

Although we look for commonalities to find labels and look for variations to find data values, we must be careful about being too strict. Sometimes there are additional or missing label-value pairs. The two nested tables beginning with *Gene Model* in Figures 1 and 3 do not share exactly the same structure. The table in Figure 1 has five columns and three rows, while the table in Figure 3 has six columns and two rows. Although they are not exactly the same, we can still identify the structure pattern by comparing them. The top rows in the two tables are very similar. Observe that the table in Figure 3 only has an additional *Swissprot* column inserted between the *Protein* and *Amino Acids* columns.

In addition to discovering the structure pattern for a web site, we can also dynamically adjust the pattern if the system encounters a table that varies from the pattern. If there is an additional or missing label, the system can change the pattern by either adding the new label and marking it optional or marking the missing label optional. For example, if we had not seen the extra *Swissprot* column in our initial pair of sibling pages, the system can add *Swissprot* as a new label and mark it as optional. The basic label-value association pattern is still the same.

We call our table-interpretation system *TISP* (*Table Interpretation with Sibling Pages*) [39].

[Home](#) [Genome](#) [Blast / Blat](#) [WormMart](#) [Batch Sequences](#) [Markers](#) [Genetic Maps](#) [Submit](#) [Searches](#) [Site Map](#)

Find: Anything [WormBase Banner](#)

[Gene Summary](#) [Locus Summary](#) [Sequence Summary](#) [Protein Summary](#) [EST Alignments](#) [Genome Browser](#) [Genetic Map](#) [Nearby Genes](#) [Tree Display](#) [Acedb XML Schema](#) [Image](#)

Gene Summary for *dyb-1*

Specify a gene using a gene name ([unc-26](#)), a predicted gene id ([R13A5.9](#)), or a protein ID ([CE02711](#)):

Identification	Location	Function	gene ontology	reactome knowledgebase	alleles	similarities	reagents	bibliography																						
Identification IDs: <table border="1" style="margin-left: 20px;"> <thead> <tr> <th>CGC name</th> <th>Sequence name</th> <th>Other name(s)</th> <th>WB Gene ID</th> <th>Version</th> </tr> </thead> <tbody> <tr> <td>dyb-1 - (<i>DYstroBrevin homolog</i>) (via person: Laurent Segalat)</td> <td>F47G6.1</td> <td>NM_058459 (inferred automatically) 1B963 (inferred automatically)</td> <td>WBGene00001115</td> <td>1</td> </tr> </tbody> </table> <p> Concise Description: The <i>dyb-1</i> gene encodes a homolog of mammalian alpha-dystrobrevin (DTNA; OMIM:601239), mutation of which can lead to left ventricular noncompaction with congenital heart defects. [details] NCBI KOGs⁺: Beta-dystrobrevin [KOG4301] Species: <i>Caenorhabditis elegans</i> NCBI: [Entrez Genes: 14670171] [AceView: 1B963] Gene model(s): <table border="1" style="margin-left: 20px;"> <thead> <tr> <th>Gene Model</th> <th>Status</th> <th>Nucleotides (coding/transcript)</th> <th>Protein</th> <th>Swissprot</th> <th>Amino Acids</th> </tr> </thead> <tbody> <tr> <td>F47G6.1 1, 2</td> <td>confirmed by cDNA(s)</td> <td>1773/7391 bp</td> <td>WP:CE26812</td> <td>DTN1 CAEEL</td> <td>590 aa</td> </tr> </tbody> </table> </p> <p> Putative ortholog(s): <i>Caenorhabditis briggsae</i>: CBG22285 [syntenic alignment] (Stein LD et al. ; best reciprocal blastp match-seg-off) </p>	CGC name	Sequence name	Other name(s)	WB Gene ID	Version	dyb-1 - (<i>DYstroBrevin homolog</i>) (via person: Laurent Segalat)	F47G6.1	NM_058459 (inferred automatically) 1B963 (inferred automatically)	WBGene00001115	1	Gene Model	Status	Nucleotides (coding/transcript)	Protein	Swissprot	Amino Acids	F47G6.1 1, 2	confirmed by cDNA(s)	1773/7391 bp	WP:CE26812	DTN1 CAEEL	590 aa	Location Genetic Position: I:-15.38 +/- 0.361 cM [mapping data] Genomic Position: I:1483084..1490474 bp	Function Mutant Phenotype: Definitions of abbreviations used in the text. RNAi Phenotype(s): WT [For details see: Ahringer JA 16 Nov 2000] WT [For details see: Rual JF 01 Sep 2004]						
CGC name	Sequence name	Other name(s)	WB Gene ID	Version																										
dyb-1 - (<i>DYstroBrevin homolog</i>) (via person: Laurent Segalat)	F47G6.1	NM_058459 (inferred automatically) 1B963 (inferred automatically)	WBGene00001115	1																										
Gene Model	Status	Nucleotides (coding/transcript)	Protein	Swissprot	Amino Acids																									
F47G6.1 1, 2	confirmed by cDNA(s)	1773/7391 bp	WP:CE26812	DTN1 CAEEL	590 aa																									

Figure 3: A Second Sample Table from WormBase.

Given that we can interpret a table, we can immediately conceptualize it and add semantic annotation to it. We augmented TISP by adding two new functions: conceptualization and annotation. We call the new system TISP++. Given a structure pattern, TISP++ can automatically generate an OWL ontology that conceptualize the pattern. TISP++ uses an OWL class to represent a table, an OWL object property to represent the nesting between two tables, and an OWL data type property to represent a label. The ontology also declares constraints such as optional and functional if applicable. After the OWL ontology is generated, TISP++ can also automatically annotate all the sibling pages with respect to this ontology. The annotated information is available to any SPARQL query platform, so that users can query and locate information of interest. By doing so, TISP++ makes hidden web information visible to users from outside specialized hand-built GUIs.

We present the details of TISP and our contribution to table interpretation by sibling page

comparison, and the details of TISP++ and our contribution to ontology generation and semantic annotation in the remainder of the paper as follows. Section 2 provides the details about how TISP analyzes a source page to recognize all HTML tables and how it decomposes nested tables, if any. Section 3 introduces the matching algorithms we use. Section 4 describes how we interpreted various matching results and find data tables. Section 5 explains how TISP infers the general structure patterns of a web site and therefore how it interprets the tables from the site. Section 5 also explains how to automatically adjust the generated patterns when variations are encountered. In Section 6, we report the results of experiments we conducted involving sites for car advertisements, molecular biology, and geopolitical information, which we found on the hidden web. Section 7 explains how TISP++ generates an OWL ontology depending on a structure pattern, and Section 8 explains how TISP++ annotates information automatically. Section 9 discusses related work. In Section 10, we draw conclusions and mention some possibilities for future work.

2 Initial Table Processing

After obtaining a source document, TISP first parses the source code and locates all HTML components enclosed by `<table>` and `</table>` tags (tagged tables). When tagged tables are nested inside of one another, TISP finds them and unnests them. In Figure 1, there are several levels of nesting in the large rectangular table. The first level is a table with two columns. The first column contains *Identification*, *Location*, and *Function*, and the second column contains some complex structures. Figure 1 shows only the first three rows of this table — one row for *Identification*, one for *Location*, and one for *Function*. (For the purpose of being explicit in this paper, we assume that these three rows are the only rows in this table.) The second column of the large rectangular table in Figure 1 contains three second level nested tables, the first starting with *IDs*, the second with *Genetic Position*, and the third with *Mutant Phenotype*. In the right most cell of the first row is another table. There are also two third level nested tables.

We treat each tagged table as an individual table and assign an identifying number to it. If the table is nested, we replace the table in the upper level with its identifying number. By so doing, we are able to remove nested tables from upper level tables. As a result, TISP decomposes the page in Figure 1 into the set of tables in Figure 4.

3 Table Matching

To compare and match tables, we first transform each HTML table into a DOM tree [18]. $Tree_1$ in Figure 5 shows the DOM tree for Table 7 in Figure 4, and $Tree_2$ in Figure 5 shows the DOM tree for its corresponding table in Figure 3.

Home	Genome	Blast / Blast	WormMart	Batch Sequences	Markers	Gen
Gene Summary	Locus Summary	Sequence Summary	Table 2			
			Protein Summary	EST Alignments	Genome Browser	Genetic M
Table 3						
[identification] [location] [function] [gene ontology] [reactome knowledgebase] [alleles] [similarities] [reagents] [bibliography]						
Table 4						
Identification		Table 5				
Location		Table 8				
Function		Table 9				
Table 5						
IDs:	Table 6					
NCBI KOGs*:	Protein kinase PCTAIRE and related kinases [KOG0594]					
Species:	<i>Caenorhabditis elegans</i>					
Other sequence(s):	AF083878 (<i>Caenorhabditis elegans</i> cyclin-dependent kinase CDK-4 (cdk-4) mRNA, complete cds.)					
NCBI:	[Entrez Genes: 15718266] [AceView : XO136]					
Gene model(s):	Table 7					
Putative ortholog(s):	<i>Caenorhabditis briggsae</i> : CBG07433 [syntenic alignment] (Stein LD et al best reciprocal blastp match-seg-off)					
Table 6						
<u>CGC name</u>	<u>Sequence name</u>	<u>Other name(s)</u>	<u>WB Gene ID</u>	<u>Version</u>		
cdk-4 - (<i>Cyclin-Dependent Kinase family</i>) (via person: Michael Krause)	F18H3.5	NM_077855 (inferred automatically) XO136 (inferred automatically)	WBGene00000406	1		
Table 7						
<u>Gene Model</u>	<u>Status</u>	<u>Nucleotides (coding/transcript)</u>	<u>Protein</u>	<u>Amino Acids</u>		
F18H3.5a 1, 2	confirmed by cDNA(s)	1029/3051 bp	WP-CE18608	342 aa		
F18H3.5b 1, 2, 3	partially confirmed by cDNA(s)	1221/1704 bp	WP-CE28918	406 aa		

Figure 4: Decomposition for the Tables in Figure 1.

Tai [38] gives a well acknowledged formal definition of the concept of a tree mapping for labeled ordered rooted trees.

Let T be a labeled ordered rooted tree and let $T[i]$ be the i^{th} node in level order of tree T . A *mapping* from tree T to tree T' is defined as a triple (M, T, T') , where M is a set of ordered pairs (i, j) , where i is from T and j is from T' , satisfying the following conditions for all $(i_1, j_1), (i_2, j_2) \in M$, where i_1 and i_2 are two nodes from T and j_1 and j_2 are two nodes from T' :

- (1) $i_1 = i_2$ iff $j_1 = j_2$;
- (2) $T[i_1]$ comes before $T[i_2]$ iff $T'[j_1]$ comes before $T'[j_2]$ in level order;
- (3) $T[i_1]$ is an ancestor of $T[i_2]$ iff $T'[j_1]$ is an ancestor of $T'[j_2]$.

According to this definition, each node appears at most once in a mapping — the order between sibling nodes and the hierarchical relation between nodes being preserved. The best match between two trees is a mapping with the maximum number of ordered pairs.

We use a simple tree matching algorithm introduced in [48] which was first proposed to compare two computer programs in software engineering. It calculates the similarity of two trees by finding the best match through dynamic programming with complexity $O(n_1 n_2)$, where n_1 is the size (number of nodes) of T and n_2 is the size of T' . This algorithm counts the matches of all possible combination pairs of nodes from the same level, one from each tree, and finds the pairs with maximum matches. The simple tree match algorithm returns the number of these maximum matched pairs. The highlighted part in $tree_1$ in Figure 5 shows the matched nodes for $tree_1$ with respect to $tree_2$ in Figure 5. The highlighted nodes indicate a match.

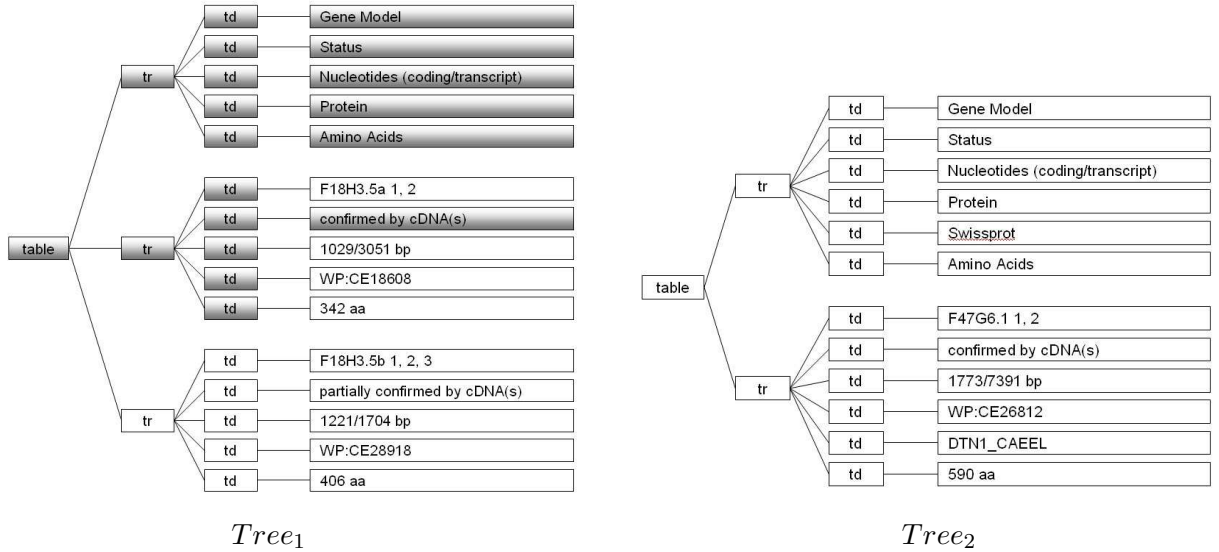


Figure 5: DOM Trees for Table 7 in Figure 4 and its Sibling Table in Figure 3.

4 Sibling Table Identification

In our research, we use the results of the simple tree matching algorithm for three tasks: (1) we filter out those HTML tables that are only for layout; (2) we identify the corresponding tables (sibling tables) from sibling pages; and (3) we match nodes in a sibling table pair.

For each pair of trees, we use the simple tree matching algorithm to find the maximum number of matched nodes among the two trees. We call this number the *match score*. For each table in one source page, we obtain match scores. Sibling tables should have a one to one correspondence. Based on the match score, we use the Gale-Shapley stable marriage algorithm [23] to pair potential sibling tables one to one.

For a pair of potential sibling tables, we calculate the *sibling table match percentage*, 100 times the match score divided by the number of nodes of the smaller tree. The match percentage between the two trees in Figure 5, for example, is 19 (match score) divided by 27 (tree size of $Tree_2$), which, expressed as a percentage, is 70.4%.

We classify potential sibling tables into three categories: (1) exact match or near exact match; (2) false match; and (3) sibling-table match. We use two threshold boundaries to classify potential sibling tables: a higher threshold between exact or near exact match and sibling-table match, and a lower threshold between sibling-table match and false match. Usually a large gap exists between the range of exact or near exact match percentages and the range of sibling-table match percentages, as well as between the range of sibling-table match percentages and the range of false match percentages. Using active learning with bootstrap selective sampling [42], we first set initial thresholds by empirical observation (90% for the higher threshold and 20% for the lower

threshold); then TISP dynamically adjusts the two thresholds as needed during the classification process as more sibling pages are considered.

In our example, Tables 1, 2, and 3 have match percentages of 100% with their sibling tables. The match percentages for Tables 4, 5, 6 and 7, and their corresponding sibling tables, are 66.7%, 58.8%, 69.2%, and 70.4% respectively. Our example has no false matches. A false match usually happens when a table does not have a corresponding table in the sibling page. In this case, we save the table. When more sibling pages are compared, we might find a matching table.

5 Structure Patterns

The first component of a structure pattern for a table specifies the table's location in a web page. To specify the location, we use XPath [2], which describes the path of the table from the root HTML tag of the document. For example, The location for Table 7 in Figure 4 is: `/html/table[4]/tbody/tr[1]/td[2]/table[2]/tbody/tr[1]/td[2]`. An XPath simply lists the nodes (HTML tag names) of a path in a DOM tree for the HTML document where [n] designates the nth sibling node in the ordered subtree.

The second component of a structure pattern specifies the label-value pairs for a table and thus provides the interpretation. We now give the details about how we identify the proper label-value pattern template (Section 5.1) and use it to generate the specific label-value-pair pattern for the table (Section 5.2). We then explain how TISP uses the generated pattern to extract label-value pairs from the table and how TISP produces an interpretation for the table (Section 5.3). Combinations of basic patterns are also possible; we thus also explain how to generate and use combination patterns (Section 5.4). Finally, we explain how TISP dynamically adjusts a pattern to accommodate table variations it may encounter as it extracts label-value pairs from sibling tables in the web site (Section 5.5).

5.1 Pattern Templates

We use regular expression to describe table structure pattern templates. If we traverse a DOM tree, which is ordered and labeled, in a preorder traversal, we can layout the tree labels textually and linearly. We can then use regular-expression like notation to represent the table structure patterns (see Figure 6). In both templates and generated patterns we use standard notation: ? (optional), + (one or more repetitions), and | (alternative). In templates, we augment the notation as follows. A variable (e.g. n) or an expression (e.g. $n-1$) can replace a symbol to designate a specific number of repetitions, which is unknown but fixed for the expression as it is applied. A pair of braces { } indicates a leaf node. A capital letter L is a position holder for a label and a capital letter V is a position holder for value. The part in a box is an *atomic pattern* which we use for combinational structural patterns in Section 5.4.

Pattern 1:
 $\langle table \rangle \langle tbody \rangle ? \left(\langle tr \rangle \left(\langle td|th \rangle \{L\}^n \mid \langle tr \rangle \left(\langle td|th \rangle \{V\}^n \right) \right)^+ \right)$
 Pattern 2:
 $\langle table \rangle \langle tbody \rangle ? \left(\langle tr \rangle \left(\langle td|th \rangle \{L\} \langle td|th \rangle \{V\}^n \right)^+ \right)$
 Pattern 3:
 $\langle table \rangle \langle tbody \rangle ? \left(\langle tr \rangle \left(\langle td|th \rangle \{L\}^n \right) \left(\langle tr \rangle \left(\langle td|th \rangle \{L\} \langle td|th \rangle \{V\}^{(n-1)} \right) \right)^+ \right)$

Figure 6: Some Basic Pre-defined Pattern Templates.

Figure 6 shows three basic pre-defined pattern templates. Pattern 1 is for tables with n labels in the first row and with n values in each of the rest of the rows. The association between labels and values is column-wise; the label at the top of the column is the label for all the values in each column. Pattern 2 is for tables with labels in the left-most column and values in the rest of the columns. Each row has a label followed by n values. The label-value association is row-wise; each label labels all values in the row. Pattern 3 is for two-dimensional tables with labels on both the top and the left. Each value in this kind of table associates with both the row header label and the column header label. We can define additional patterns, but these patterns and combinations of these patterns constitute a large majority of HTML tables.

5.2 Pattern Generation

To check whether a table matches any pre-defined pattern template, TISP tests each template until it finds a match. When we search for a matching template, we only consider leaf nodes and seek matches for labels and mismatches for values. Variations, however, exist and we must allow for them. In tables, labels or values are usually grouped. We are seeking for a structure pattern instead of classifying individual cells. Sometimes we find a matched node, but all other nodes in the group are mismatched nodes and agree with a certain pattern, TISP should ignore the disagreement and assume the matched node is a mismatched node of values too. Specifically, we calculate a *template match percentage* between a pre-defined pattern template and a matched result, 100 times the number of leaf nodes that agree with a pattern template divided by total number of leaf nodes in the tree. We calculate the template match percentage between a table and each pre-defined structure template. A match must satisfy two conditions: (1) it must be the highest match percentage, and (2) the match percentage must be greater than a threshold. Similar to the way we determine thresholds for sibling table matches, we determine this template match percentage threshold using active learning with bootstrap selective sampling, with an initial threshold of 80%.

Consider the mapped result in Figure 5 as an example. The highlighted nodes are matched nodes in $tree_1$. Comparing the template match percentage for this mapped result for the three pattern templates in Figure 6, we obtain 93.3%, 53.3%, and 80% respectively. Pattern 1 has the

```

/html/table[4]/tbody/tr[1]/td[2]/table[2]/tbody/tr[1]/td[2]
< table >< tr >
< td > Gene Model
< td > Status
< td > Nucleotides(coding/transcript)
< td > Protein
< td > Amino Acids
(< tr >
< td > VGene Model
< td > VStatus
< td > VNucleotides(coding/transcript)
< td > VProtein
< td > VAmino Acids)+

```

Figure 7: Structure Pattern for Table 7 in Figure 4.

highest match percentage, and it is greater than the threshold. Therefore we choose Pattern 1.

We now impose the chosen pattern, ignoring matches and mismatches. Note that for $tree_1$ in Figure 5, the first branch matches the part in Pattern 1 in the first box, and the second and the third branch each match the part in the second box, where n is five. For Pattern 1, when $n=1$, we have a one-dimensional table; and when $n>1$, we have a two-dimensional table for which we must generate record numbers.

After TISP matches a table with a pre-defined pattern template, it generates a specific structure pattern for the table by substituting the actual labels for each L and by substituting a placeholder V_L for each value. The subscript L for a value V designates the label for the label-value pair for each record in a table. Figure 7 shows the specific structure pattern for Table 7 in Figure 4.

5.3 Pattern Usage

With a structure pattern for a specific table, we can interpret the table and all its sibling tables. The XPath gives the location of the table, and the generated pattern gives the label-value pairs. The pattern must match exactly in the sense that each label string encountered must be identical to the pattern's corresponding label string. Any failure is reported to TISP. (In Section 5.5, we explain how TISP reacts to a failure notification.)

When the pattern matches exactly, TISP can generate an interpretation for the table. For our example, the chosen pattern is Pattern 1 with $+$ (which allows for multiple rows of values in the table). Thus, TISP needs to add another dimension and add row numbers. Since the table is inside of other tables, TISP recursively searches for the tables in the upper levels of nesting and collects all needed labels.

This process repeats recursively until all sub-groups match with a template or the process fails to finding any matching template.

For the example in Figure 8, TISP is unable to find any rows of labels, but finds two columns of labels, the first and third column. It then divides the table into two groups using these two columns and tries to match each group with a pre-defined template. It matches each group with Pattern 2. Therefore, this table matches column-wise with Pattern 2 used twice.

5.5 Dynamic Pattern Adjustment

Given a structure pattern for a table, we know where the table is in the source document (its XPath), the location of the labels and values, and the association between labels and values. When TISP encounters a new sibling page, it tries to locate each sibling table following the XPath, and then tries to interpret it by matching it with the sibling table structure pattern. If the encountered table matches the structure pattern regular expression perfectly, we successfully interpret this table. Otherwise, we might need to do some pattern adjustment. There are two ways to adjust a structure pattern: (1) adjust the XPath to locate a table, and (2) adjust the generated structure pattern regular expression.

Although sibling pages usually have the same base structure, some variations might exist. Some sibling pages might have additional or missing tables. Thus, sometimes, following the XPath, we cannot locate the sibling table for which we are looking. In this case, TISP searches for tables at the same level of nesting, looking for one that matches the pattern. If TISP finds one, it obtains the XPath and adds it as an alternative. Thus, for future sibling pages, TISP can (in fact, always does) check all alternative XPaths before searching for another alternative XPath. If TISP finds no matching table, it simply continues its processing with the next table.

We adjust a table pattern when we encounter a variation of an existing table. There might be additional or missing labels in the encountered variation. In this case, we need to adjust the structure pattern regular expression, to add the new optional label or to mark the missing label as optional. Consider the table that starts with *Gene Model* in Figure 3 (the sibling table of Table 7 in Figure 4) as an example. The table matches the pattern in Figure 7 until we encounter the label *Swissprot*. If we skip *Swissprot*, the next label *Amino Acids* matches the structure pattern. In this case, we treat *Swissprot* as an additional label, and we add it as an optional label as Figure 10 shows.

6 Experimental Results

We tested TISP using source pages from commercial data, scientific data, and geopolitical data. We picked pages from each field: car advertisements for commercial data, molecular biology for scientific data, and interesting information about US states and about countries for geopolitical

```

< table >< tr >< td >Gene Model< td >Status < td >Nucleotides(coding/transcript)
< td >Protein (< td >Swissprot)? < td >Amino Acids
(< tr > < td >VGene Model< td >VStatus < td >VNucleotides(coding/transcript)
< td >VProtein(< td >VSwissprot)? < td >VAmino Acids)+

```

Figure 10: Structure Pattern for the Table in Figure 3.










data. Most of the source pages were collected from popular and well-known web sites such as cars.com, NCBI database, Wormbase, MTB (Mouse Tumor Biology Database), CIA's World Factbook, and U.S. Geological Survey (usgs.gov). We tested more than 2,000 tables found in 275 sibling pages in 35 web sites. Most pages from the molecular biology domain and the geopolitical domain have relatively complicated structures. Seven out of ten sites in the geopolitical domain and eight out of ten sites in the molecular biology domain contained multiple data tables per page. Two of the geopolitical sites and eight of the molecular biology sites contained nested HTML tables.

For each web site, we randomly chose two sibling pages for initial pattern generation. For the initial two sibling pages, we tested (1) whether TISP was able to recognize HTML data tables and discard HTML tables used only for layout, (2) whether it was able to pair all sibling tables correctly, and (3) whether it was able to recognize the correct pattern template or pattern combination. For the rest of sibling pages from the same web site, we tested (1) whether TISP was able to interpret tables using the recognized structure patterns, (2) whether it correctly detected the need for dynamic adjustment, and (3) whether it recognized new structure patterns correctly.

We collected 75 sibling pages from 15 different web sites in the car-advertisements domain for a total of 780 HTML tables.³ TISP correctly discarded all uses of tables for layout and successfully paired all sibling tables. There were no nested tables in this domain. Most of the web sites contained only one table pattern, except for one site that had three different patterns. Two web sites contained tables with structure combinations. Of the 17 pairs of sibling tables, TISP recognized 16 correctly. The one pattern TISP failed to recognize correctly contained too many value cells that included the same value (values such as *unknown*, *general car*, *auto*, *dealer*, and empty spaces). TISP considered them as labels, and thus could not match the table with any pre-defined pattern template or detect any pattern combination. TISP successfully interpreted all tables from the generated patterns. No adjustment were needed, neither for any path nor for any label.

For the geopolitical information domain, we tested 100 sibling pages from 10 different web sites with 884 HTML tables. TISP correctly paired 100% of all data tables and correctly discarded all layout tables. For initial pattern generation, TISP was able to recognize all 22 structure patterns

³The sibling pages in this domain are usually very regular. Indeed, we found no table variations in any of the sites we considered. We, therefore, only tested five pages per site.

Geography		Mongolia	
Location:			
Northern Asia, between China and Russia			
Geographic coordinates:			
46 00 N, 105 00 E			
Map references:			
Asia			
Area:			
total: 1,564,116 sq km			

(a)

Crime in New York by Year

Type	1999
Murders	671
per 100,000	8.4
Rapes	1,702
per 100,000	21.3

(b)

Figure 11: Two Partially Misinterpreted Tables.

successfully. As TISP processed additional sibling pages, it found one additional sibling table and correctly interpreted it. There were no path adjustments, but there were 22 label adjustments — all of them correct. For two sets of sibling tables, TISP recognized the correct patterns, but failed to recognize some implicit information that affects the meaning of the tables. Therefore it interpreted these tables only partially correctly. Figure 11 shows these two cases. There are actually two HTML tables in Figure 11a. The header *Geography Mongolia* is in one HTML table, and the rest of information is in another HTML table. Because it separated tables using HTML tags, TISP was not able to determine the relationship between these two HTML tables. TISP correctly interpreted Figure 11b as Pattern 3. It, however, did not recognize the relationship between *Murders* and *per 100,000* and between *Rapes* and *per 100,000*.

We collected 100 sibling pages from 10 different web sites in the molecular biology domain for a total of 862 HTML tables. Among these tables, TISP falsely classified three pairs of layout tables as data tables. TISP, however, successfully eliminated these false sibling pairs during pattern generation because it was unable to find a matching pattern. No false patterns were generated. TISP was able to recognize 28 of 29 structure patterns. TISP missed one pattern because the table contained too many empty cells. If it had considered empty cells as mismatches, TISP would have correctly recognize this pattern. As TISP processed additional sibling pages, it found 5 additional sibling tables and correctly interpreted all but one of them. The failure was caused by

labels that varied across sibling tables causing them, in some cases, to look like values. There were 5 path adjustments and 12 label adjustments — all of them correct. One table was interpreted only partially correctly because TISP considered the irrelevant information *To Top* as a header.

For measuring the overall accuracy of TISP, we computed precision (P), recall (R), and an F-measure ($F = 2PR/(P+R)$). In its table recognition step, TISP correctly discarded 155 of 158 layout tables and discarded no data tables. It therefore achieved an F-measure of 99.0% (98.1% recall and 100% precision). TISP later discarded these three layout tables in its pattern generation step, but it also rejected two data tables, being unable to find any pattern for them. It thus achieved an F-measure of 99.4% (100% recall and 98.8% precision). For table interpretation, TISP correctly recognized 69 of 74 structure patterns. It therefore achieved a recall of 93.2%. Of the 72 structure patterns it detected, 69 were correct. It therefore achieved a precision of 95.8%. Overall the F-measure for table interpretation was 94.5% for the sites we tested.

We discuss the time performance of TISP in two phases: (1) initial pattern generation from a pair of sibling pages and (2) interpretation of the tables in the rest of the sibling pages. The time for the pattern generation given a pair of sibling pages consists of: (1) the time to read and parse the two pages and locate all the HTML tables, (2) the time for sibling table comparisons, and (3) the time to select from pre-defined structure templates and generate a pattern. The complexity of parsing and locating HTML tables is $O(n)$, where n is the number of HTML tags. The simple tree matching algorithm has time complexity $O(m_1m_2)$, where m_1 and m_2 are the numbers of nodes of the two sibling trees. To find the best match for each HTML table, we need to compare each table with all the HTML tables in its sibling page. The time complexity is $O(km_1m_2)$, where k is the number of HTML tables in the sibling page. The time complexity for finding the correct pattern for each matched sibling table is $O(pl)$, where p is the number of pattern templates and l is the number of leaf nodes in the HTML table. If pattern combinations are involved, the complexity of template discovery increases multiplicatively since for each subgroup we must consider every template and find the best match. We conducted the experiment on a Pentium 4 computer running at 3.2 GHz. The typical actual time needed for the pattern generation for a pair of sibling pages was below or about one second. The actual time reached a maximum of 15 seconds for a complicated web site where pages had more than 20 tables.

The time for table interpretation for a single sibling web page when no adjustment is necessary consists of: (1) the time for locating each table and (2) the time for processing the table with a pattern. The complexity of locating a table is $O(p)$, where p is the number of path possibilities leading to the table. Each path possibility is itself logarithmic with respect to the number of nodes in the DOM tree for the pages. The complexity of matching a located table with the corresponding pattern is $O(el)$, where e is the number of pattern entries (an entry could be either a pattern label or a pattern value) of the pattern and l is the number of leaf nodes in the HTML table's DOM

tree. The time to do adjustments ranges from the time to do a simple label adjustment, which is constant, to the time required to re-evaluate all sibling tables, which is the same as the time for initial pattern generation. Overall, the typical actual time needed for interpreting tables in one page was below one second. The actual time reached a maximum of 19 seconds for a complicated web page with several tables and several adjustments.

7 Semantic Ontology Generation

After we obtain the structure pattern for a web site, we can conceptualize the information present in the tables by generating an ontology according to the structure pattern. TISP++ uses a structure pattern to generate OWL classes, properties, and constraints according to the structure pattern. It then uses Jena [3], a semantic web framework for Java, to output the OWL ontology.

For each web site, we choose a name. This name provides an anchor to which we attach ontological concepts. For example, we call the WormBase gene repository “WormBaseGene”. TISP++ generates the class “WormBaseGene” for the Worm Base gene repository. Line 5 in Figure 12 shows the OWL class “WormBaseGene”.

For each label, TISP++ either generates an OWL class or a data type property. If a label is paired with a nested table such as are the labels *Identification* and *Gene model(s)* in Figures 1 and 3, TISP++ generates an OWL class for the label as Lines 9 and 11 in Figure 12 show. If a label is paired with an actual value such as are the labels *Gene Model* and *Amino Acids* in Figure 1 and 3, TISP++ generates an OWL data type property for this label, as Lines 25, 29, and 34 in Figure 12 show.

The generated ontology also represents the relationships among the labels. TISP++ interprets tables according to the structure patterns in Figure 6. For patterns such as Pattern 1 and Pattern 2 that only involve binary relationships, relationship generation is straightforward. For the relationship between a label *A* defined by an OWL class and a label *B* defined by an OWL data type property (such as the relationship between *Gene Model(s)* and *Gene Model* in Figure 1), TISP++ defines *A* as the domain and *B* as the range with *string* as its type by default. As Figure 12, Lines 21-24 show, we define the *DatatypeProperty GeneModel* as a *string* whose domain is the already defined class *Genemodels*. If the data type property is optional, as is *Swissprot* in Figure 10, we add a *minCardinality* declaration allowing zero occurrences. Line 26 in Figure 12 shows this *minCardinality* declaration for *Swissprot*. For the relationships between two classes *A* and *B*, TISP++ generates an OWL object property: *A--B* and its inverse *B--A*. For the property *A--B*, TISP++ defines *A* as the domain and *B* as the range. For example, Lines 9–15 in Figure 12 show the OWL object property for *WormBaseGene--Identification*.

For patterns such as Pattern 3 that involve n-ary relationships, relationship generation requires a preparatory transformation. Since OWL ontologies only allow us to declare binary relationships,

```

1. <rdf:RDF
2.     xmlns:owl="http://www.w3.org/2002/07/owl#"
3.     xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
4.     xmlns:base="http://dithers.cs.byu.edu/owl/ontologies/wormbase#" ... >
...
5.     <owl:Ontology rdf:about="" >
6.         <rdfs:comment>OWL Ontology for WormBase </rdfs:comment>
7.         <rdfs:label>WormBase Ontology</rdfs:label>
8.     </owl:Ontology>
9.     <owl:Class rdf:ID="WormBaseGene" />
10.    <owl:Class rdf:ID="Identification" />
11.    <owl:Class rdf:ID="Genemodels" />
12.    ...
13.    <owl:ObjectProperty rdf:ID="WormBaseGene- -Identification" >
14.        <owl:inverseOf>
15.            <owl:ObjectProperty rdf:ID="Identification- -WormBaseGene" >
16.        </owl:inverseOf>
17.        <rdfs:range rdf:resource="#Identification" />
18.        <rdfs:domain rdf:resource="#WormBaseGene" />
19.    </owl:ObjectProperty>
20.    ...
21.    <owl:ObjectProperty rdf:ID="Identification- -IDs" >
22.    ...
23.    <owl:ObjectProperty rdf:ID="Identification- -Genemodels" >
24.    ...
25.    <owl:DatatypeProperty rdf:ID="GeneModel" >
26.        <rdfs:range rdf:resource="http://www.w3.org/2001/XMLSchema#string" />
27.        <rdfs:domain rdf:resource="#Genemodels" />
28.    </owl:DatatypeProperty>
29.    ...
30.    <owl:DatatypeProperty rdf:ID="Swissprot" >
31.        <owl:minCardinality>0</owl:minCardinality>
32.        <rdfs:range rdf:resource="http://www.w3.org/2001/XMLSchema#string" />
33.        <rdfs:domain rdf:resource="#Genemodels" />
34.    </owl:DatatypeProperty>
35.    <owl:DatatypeProperty rdf:ID="AminoAcids" >
36.        <rdfs:range rdf:resource="http://www.w3.org/2001/XMLSchema#string" />
37.        <rdfs:domain rdf:resource="#Genemodels" />
38.    </owl:DatatypeProperty>
39. </rdf:RDF>

```

Figure 12: Partial OWL Ontology for WormBase.

we transform Pattern 3 in Figure 6 by considering the label L in the third line “($\langle tr \rangle \langle td|th \rangle \{L\} \langle td|th \rangle \{V\}^{(n-1)+}$ ” as a value V . Then it becomes “($\langle tr \rangle \langle td|th \rangle \{V\} \langle td|th \rangle \{V\}^{(n-1)+}$ ” and can be further simplified to “($\langle tr \rangle \langle td|th \rangle \{V\}^n$)”, which is the same as Pattern 1. Therefore TISP++ can transform n-ary relationships to binary relationships and then translate them as binary relationships to OWL.

8 Semantic Annotation and Querying

After TISP++ generates an ontology according to the structure pattern of a web repository, it automatically annotates the pages from this repository with respect to the generated ontology. Figure 13 shows a partial graphical view of the annotated RDF triples for the page in Figure 1. An oval with a solid border represents an OWL class. A dashed oval represents a URI instance. A rounded rectangle represents a value. A line represents an OWL property, and the text on the line indicates the name space and the name of the property. For example, the URI instance *#WormBaseGene*, which refers to the whole table, is a URI instance for the *WormBaseGene* class in the *wormbase* ontology in Figure 12. The URI instance *#Identification*, which refers to the value of the label *Identification* in Table 4 in Figure 4, is a URI instance of the class *Identification* in the *wormbase* ontology. Since Table 7 in Figure 4 has two data rows, TISP++ declares URI instances: *#Genemodels1* and *#Genemodels2*, one for each row. TISP++ also declares the relationship between two URI instances. For example, TISP++ declares relationships *WormBaseGene-Identification* and its reverse *Identification-WormBaseGene*. TISP++ also declares relationships between instances and values. For example, the *Protein* value for *#Genemodels1* is “WP:CE18608”.

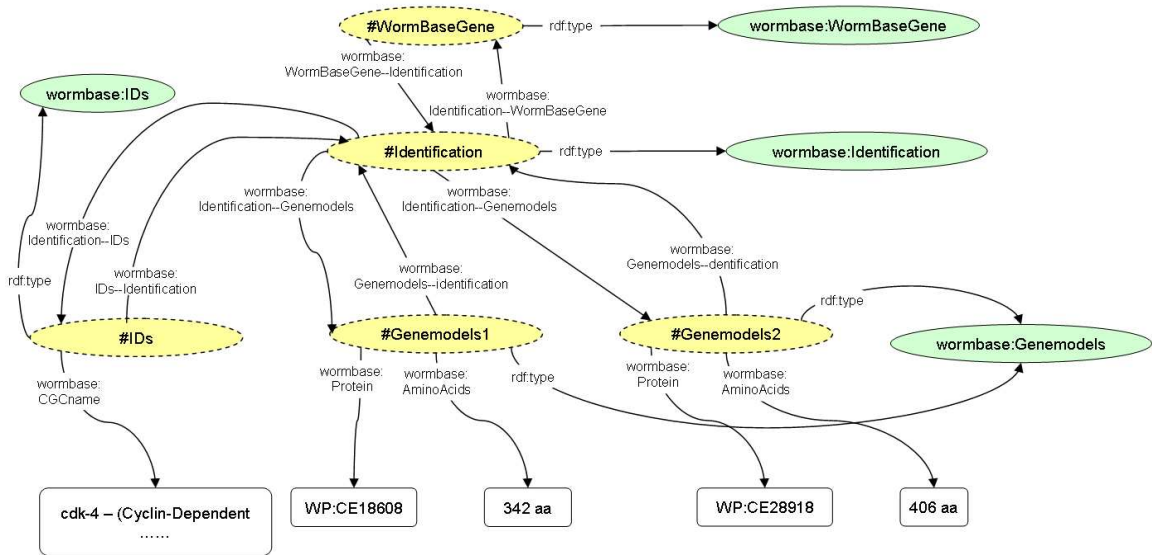


Figure 13: Graphical View of the RDF Annotation of Figure 1 for the Ontology in Figure 12.

```

1. <rdf:RDF ...
2.     xmlns:wormbase="http://www.deg.byu.edu/ontology/wormbase#"
3.     xmlns:ann="http://www.deg.byu.edu/ontology/annotationSpecification#"
4.     xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#" >
...
5.     <wormbase:WormBaseGene rdf:about="http://localhost/wormbase.html#WormBaseGene" >
6.         <wormbase:WormBaseGene- -Identification
7.             rdf:resource="http://localhost/wormbase.html#Identification" >
8.     </wormbase:WormBaseGene>
9.     <wormbase:Identification rdf:about="http://localhost/wormbase.html#Identification" >
10.        <wormbase:Identification- -WormBase> ...
11.        <wormbase:Identification- -Genemodels>
12.            <wormbase:Genemodels rdf:about="http://localhost/wormbase.html#Genemodels1" >
13.                ...
14.                <wormbase:Protein>WP:CE18608</wormbase:Protein>
15.                <wormbase:AminoAcids>342 aa</wormbase:AminoAcids>
16.                <ann:Path>/html/table[4]/tbody/tr[1]/td[2]/table[2]/tbody/
17.                    tr[1]/td[2]/table/tr[2]</wormbase:Path>
18.            <wormbase:Genemodels rdf:about="http://localhost/wormbase.html#Genemodels2" >
19.                ...
20.                <wormbase:Protein>WP:CE28918</wormbase:Protein>
21.                <wormbase:AminoAcids>406 aa</wormbase:AminoAcids>
22.                ...
23.            ...
24.        </wormbase:Identification- -Genemodels>
25.    </wormbase:Identification >
26. </rdf:RDF>

```

Figure 14: Annotation for Figure 1 for the Ontology in Figure 12.

Figure 14 shows a portion of the corresponding RDF file for the RDF graph in Figure 13. In an RDF annotation file, we first declare the name spaces of referenced ontologies. Line 2 in Figure 14 refers to the *wormbase* ontology in Figure 12, and Lines 3 and 4 declare additional name spaces. Annotated values appear as themselves, tagged with their appropriate labels. For example, the protein value “WP:CE18608” is in Line 12. To identify the annotated string in the original page, TISP++ keeps track of the original position where the values are located by recording the XPath to the node in the DOM tree of the cached page, a local copy of the original page, containing the values. Line 14 in Figure 14, for example, shows where we can locate the first row of values for Table 7 in Figure 4.

With the annotated data properly stored in an RDF file (e.g. Figure 12), we are now ready to query the annotated data using SPARQL [4]. We adapt Twinkle [5] as our SPARQL query interface. A simple query illustrates how it works. Suppose we want to find the protein and the amino-acids information for gene “cdk-4”, then we can write the SPARQL query in Figure 15. Figure 15, which is a screen shot of our adaption of Twinkle, also shows the query and the result of this query. The SPARQL query first finds the URI instance *#IDs* whose *wormbase:CGCname* contains “cdk-4” as specified in the first and last two lines of the WHERE clause. It then finds the URI instance *#Identification* through the *wormbase:IDs- -Identification* relationship and then URI instances *#Genemodels1* and *#Genemodels2* through the *wormbase:Identification- -Genemodels*

relationship. Finally it can locate the results by finding the *Protein* and *AminoAcids* values respectively through the relationships *wormbase:Protein* and *wormbase:AminoAcids*.

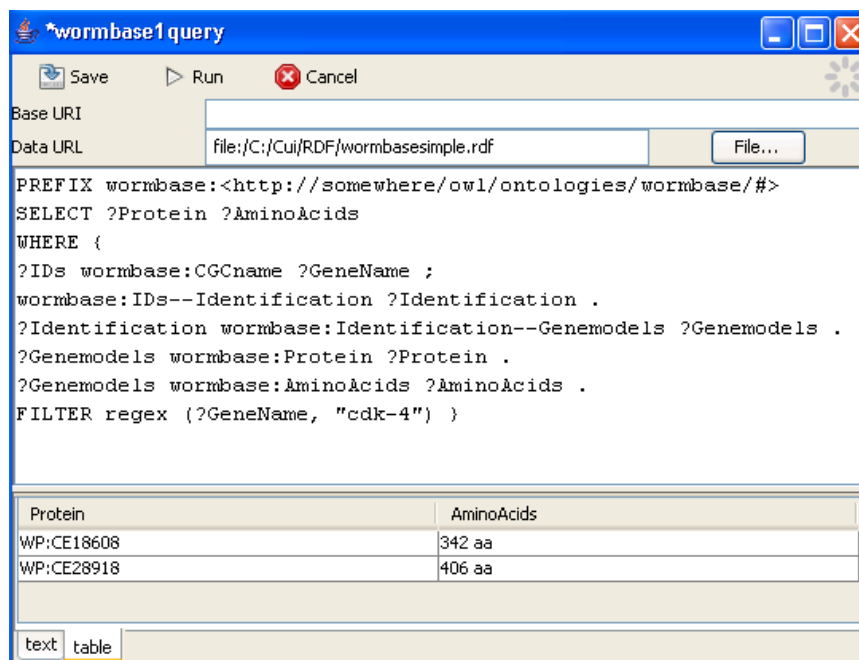


Figure 15: A Sample SPARQL Query and the Results for the Annotation in Figure 13

We now make hidden web data present in HTML tables completely accessible by a standard query system. In addition, it also semantically annotates the data with respect to the generated OWL ontology, so that the data becomes machine-understandable and automatically manipulatable by computer agents.

9 Related Work

9.1 Sibling Page Comparison

Other researchers have also tried to take advantage of sibling pages. RoadRunner [13] compares two HTML pages from one web site and analyzes the similarities and dissimilarities between them in order to generate extraction wrappers. It discovers data fields by string mismatches and discovers iterators and optionals by tag mismatches. EXALG [7] uses equivalence classes (sets of items that occur with the same frequency in sibling pages) and differentiating roles to generate extraction templates for the sibling pages. DEPTA [50] compares different records in a page instead of sibling pages and tries to find the extraction template for the record. This approach first tries to find individual data records by using a few heuristics. It then uses a tree edit distance algorithm to compare different data records and tries to find the extraction region. The approach in [30] compares sibling pages to filter out general headers and footers and other constant non-data

areas of a page. It then makes various comparisons among main pages and linked pages to find record segmentations.

TISP fundamentally differs from these approaches. The first three [7, 13, 50] focus on finding data fields, and the technique in [30] focuses on record segmentation. They do not discover labels or try to associate data and labels. TISP focuses on table interpretation. It looks for a table pattern in addition to data fields. Furthermore, TISP also tries to find the general structure pattern for the entire web site. It dynamically adjusts the structure pattern as it encounters new, yet-unseen structures.

9.2 Table Interpretation

Automated table processing is typically done in two steps: (1) table recognition — find the data table, and (2) table interpretation — find and associate labels and values. Recent surveys [20, 49] describe the vast amount of research that has been done in table processing and illustrate the challenges of automated table processing. Most of this work is about tables in imaged documents, but some is about HTML tables. Since we focus in this paper only on HTML tables, we limit the related work we discuss to only HTML table processing.

Several researches have tried to differentiate data tables from tables for layout [9, 12, 24, 46]. They have tried to use machine learning methods [12, 46], visual level features [24, 25], and linguistic features [9]. TISP provides a unique way to do this task when sibling pages are available. By considering the match percentage among potential sibling tables, we were able to filter out all the layout tables in the car and geopolitical domain and only failed to filter out three pairs of tables out of more than 800 HTML tables from the molecular biology domain (These three false positives were also filtered out during the process of pattern generation). The approaches [9, 12, 24, 46] were able to achieve F-measures of 86.5% [9], 95.5% [12], 90.0% [24], and 87.6% [46]. By way of comparison, TISP was able to achieve an F-measure of 99.4%. TISP techniques, of course, only work when sibling pages and tables are available.

Several papers have discussed the HTML table interpretation problem. Some table interpretation systems work based on simple assumptions and heuristics (e.g. [9, 21, 22, 27, 31]). These simple assumptions (labels are either the first few rows or the first few columns) are easily broken in complex tables such as nested tables (e.g. Figure 1) or tables with combination structures (e.g. Figure 8). The approach in [35] presents a table interpretation system for automatic generation of F-logic frames for tables. It considers many linguistic features in a table such as alphabetic features, numeric features, number ranges, and data formats. It calculates differences among different regions of a table to detect the orientation of a table and to locate label cells and value cells. The average F-measure of this approach is around 50%. The technique in [41] learns lexical variants from training examples and uses a vector space model to deal with non-exact matches

among labels. It also uses a few heuristics to find the association among labels and values. It achieves an F-measure of 91.4%. The approach in [25] uses visual boxes instead of HTML tags to interpret HTML tables. It achieves an F-measure of 52.1% (the precision value was 57% and the recall value was 48%). By way of comparison, TISP is able to achieve an F-measure of 94.5%. Of course, TISP techniques only work when HTML sibling tables are available. On the other hand, when applicable, TISP has the advantage over machine learning because it is unsupervised and document and web-site independent. TISP has no need for training data and works for all domains and web sites where sibling pages with sibling tables are available.

9.3 Ontology Generation

In recent years, many researches have tried to facilitate ontology generation. Manual editing tools such as Protégé [33] and OntoWeb [37] have been developed to help users create and edit ontologies. It is not trivial, however, to learn ontology modeling languages and complex tools in order to manually create ontological description for information repositories.

Because of the difficulties involved in manual creation, researchers have developed semi-automatic ontology generation tools. Most efforts so far have been devoted to automatic generation of ontologies from text files. Tools such as OntoLT [8], Text2Onto [11], OntoLearn [32], and KASO [47] use machine learning methods to generate an ontology from arbitrary text files. These tools usually require a large training corpus and use various natural language processing algorithms to derive features to learn ontologies. The results, however, are not very satisfactory [34].

Tools such as SIH [40], TANGO [43], and the one developed by Pivk [34] use structured information (HTML tables) as a source for learning ontologies. Structured information makes it easier to interpret new items and relations. The approach in [34] tries to discover semantic labels for table regions and generate an ontology based on a table's structure. But how this process is done and what format the generated ontologies have is not discussed in the paper. SIH [40] and TANGO [43] are two ongoing projects we are currently working on. SIH tries to generate user-specified ontologies depending on user-generated forms. TANGO generates ontologies by analyzing related tables in a specific domain, generating an ontology according to each table, and then merging these ontologies to a general ontology for the domain. TISP++ can generate OWL ontologies fully automatically. It, however, only generates an ontology for a single set of sibling pages. It does not merge ontologies generated from different web sites, nor does it provide for user-specified ontologies. In addition, TISP++ generates ontologies in only one simple way, while TANGO aims at generating more sophisticated ontologies.

9.4 Semantic Annotation

Existing semantic annotation systems can be classified into pattern-based systems and machine learning-based systems. Pattern-based systems such PANKOW [10] and Armadillo [17] find entities by discovering patterns. The pattern are either discovered manually or induced semi-automatically with a set of initial manually tagged seed patterns. Systems such as SemTag [14], AeroDAML [29], and KIM [36] use a set of pre-defined rules to locate the information of interest. OWL-AA [15, 16] uses a domain-specified extraction ontology to locate semantic entities. Systems such as S-CREAM [26] and MnM [44] use machine learning algorithms and natural language processing methods to locate semantic entities. All of these approaches require some pre-defined information. Pattern-based approaches need a set of initial seed patterns. Rule-based approaches need a set of pre-defined rules. Extraction-ontology-based approaches need domain ontologies. And machine learning-based approaches need a training corpus. TISP++, however, does not require predefined information of any kind and still works well when sibling pages are available.

10 Conclusion and Future Work

In this paper we introduced TISP, an approach to automatically interpret tables in hidden-web pages — pages which are almost always sibling pages. By comparing data tables in sibling pages, TISP is able to find the location of table labels and data entries and pair them to infer the general pattern for all sibling tables from the same site. Our experiments using source pages from three different domains — car advertisements, molecular biology, and geopolitical information — indicate that TISP can succeed in properly interpreting tables in sibling pages. TISP achieved an F-measure for sibling table interpretation of 94.5%.

We also extended TISP to TISP++. TISP++ uses TISP results to semantically annotate web pages, turning their embedded facts into externally accessible facts. Given an interpreted table, TISP++ automatically generates an OWL ontology depending on the table’s structure and then semantically annotates the data in the table with respect to this generated ontology. By doing so, all the data present in the sibling tables becomes accessible through a standard query interface.

Several directions remain to be pursued. For TISP and table interpretation, we would like to do the following. (1) We assumed that information in one table cell is either a table label or a table value. There could be structured information within a cell, however, such as the label *via person* and the value *Michael Krause* in Figure 1. As a future work on table interpretation, we would like to analyze cell content to find structured information within cells. (2) Some web pages use lists as tables; we would like to consider them too. (3) We would also like to interpret non-HTML tables such as tables in plain text. (4) We would also like to be able to deal with the case in Figure 11a, where we need to join adjacent HTML tables to form a single table, and we

would like to improve TISP so that it can interpret tables with factored labels such as those in Figure 11b.

For TISP++ and semantic ontology generation and semantic annotation, we would like to do the following. (1) Average web users need a more user-friendly query system, so that they can find data of interest without knowing SPARQL. We plan to provide our users with a natural-language-based query interface (e.g. like [6]) and a form-based query interface (e.g. like [19]). (2) We plan to generate ontologies according to a user's personalized view. Our users would have the option to choose which data in a table they want include in the generated ontology and how this data should be organized. (3) We would like to generate more sophisticated ontologies that cover more complicated situations such as n-ary relationships, generalization/specialization, and aggregation.

References

- [1] Worm base! <http://www.wormbase.org>, 2005.
- [2] XML Path Language (XPath). <http://www.w3.org/TR/xpath>, 2006.
- [3] Jena — A Semantic Web Framework for Java. <http://jena.sourceforge.net/>, 2008.
- [4] SPARQL Query Language for RDF. <http://www.w3.org/TR/rdf-sparql-query/>, 2008.
- [5] Twinkle: A SPARQL Query Tool. <http://www.ldodds.com/projects/twinkle/>, 2008.
- [6] M.J. Al-Muhammed. *Ontology Aware Software Service Agents: Meeting Ordinary User Needs on the Semantic Web*. PhD thesis, Brigham Young University, 2007.
- [7] A. Arasu and H. Garcia-Molina. Extracting structured data from web pages. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data (SIGMOD'03)*, pages 337–348, San Diego, California, 2003.
- [8] P. Buitelaar, D. Olejnik, and M. Sintek. Ontolt: A Protégé plug-in for ontology extraction from text based on linguistic analysis. In *Proceedings of the First European Semantic Web Symposium (ESWS'04)*, pages 31–44, Heraklion, Greece, May 2004.
- [9] H. Chen, S. Tsai, and J. Tsai. Mining tables from large scale HTML texts. In *Proceedings of the 18th International Conference on Computational Linguistics (COLING'00)*, pages 166–172, Saarbrücken, Germany, July–August 2000.
- [10] P. Cimiano, S. Handschuh, and S. Staab. Towards the self-annotating web. In *Proceedings of the 13th International Conference on World Wide Web (WWW'04)*, pages 462–471, New York, New York, May 2004.

- [11] P. Cimiano and J. Völker. Text2Onto—a framework for ontology learning and data-driven change discovery. In *Proceedings of the 10th International Conference on Applications of Natural Language to Information Systems (NLDB'05)*, pages 227–238, Alicante, Spain, June 2005.
- [12] W.W. Cohen, M. Hurst, and L.S. Jensen. A flexible learning system for wrapping tables and lists in HTML documents. In *Proceedings of the 11th International World Wide Web Conference (WWW'02)*, pages 232–241, Honolulu, Hawaii, May 2002.
- [13] V. Crescenzi, G. Mecca, and P. Merialdo. RoadRunner: Towards automatic data extraction from large web sites. In *Proceedings of the 27th International Conference on Very Large Data Bases (VLDB'01)*, pages 109–118, Rome, Italy, September 2001.
- [14] S. Dill, N. Eiron, D. Gibson, D. Gruhl, R. Guha, A. Jhingran, T. Kanungo, K.S. Mccurley, S. Rajagopalan, and A. Tomkins. A case for automated large-scale semantic annotation. *Journal of Web Semantics: Science, Services and Agents on the World Wide Web*, 1(1):115–132, December 2003.
- [15] Y. Ding, D.W. Embley, and S.W. Liddle. Automatic creation and simplified querying of semantic web content: An approach based on information-extraction ontologies. In *Proceedings of the First Asian Semantic Web Conference (ASWC'06)*, pages 400–414, Beijing, China, September 2006.
- [16] Y. Ding, D.W. Embley, and S.W. Liddle. Enriching OWL with instance recognition semantics for automated semantic annotation. In *Proceedings of the First International Workshop on Ontologies and Information Systems for the Semantic Web (ONISW'2007)*, Auckland, New Zealand, November 2007.
- [17] A. Dingli, F. Ciravegna, and Y. Wilks. Automatic semantic annotation using unsupervised information extraction and integration. In *Proceedings of the Third International Conference on Knowledge Capture (K-CAP'03), Workshop on Knowledge Markup and Semantic Annotation*, Sanibel Island, Florida, October 2003.
- [18] The W3C architecture domain. <http://www.w3.org/dom/>, 2005.
- [19] D.W. Embley. NFQL: the natural forms query language. *ACM Transactions on Database Systems*, 14(2):168–211, 1989.
- [20] D.W. Embley, M. Hurst, D. Lopresti, and G. Nagy. Table processing paradigms: A research survey. *International Journal of Document Analysis and Recognition*, 8(2-3):66–86, June 2006.

- [21] D.W. Embley, C. Tao, and S.W. Liddle. Automatically extracting ontologically specified data from HTML tables with unknown structure. In *Proceedings of the 21st International Conference on Conceptual Modeling (ER'02)*, pages 322–327, Tampere, Finland, October 2002.
- [22] D.W. Embley, C. Tao, and S.W. Liddle. Automating the extraction of data from HTML tables with unknown structure. *Data & Knowledge Engineering*, 54(1):3–28, July 2005.
- [23] D. Gale and L.S. Shapley. College admissions and the stability of marriage. *American Mathematics Monthly*, 69(1):9–14, 1962.
- [24] W. Gatterbauer and P. Bohunsky. Table extraction using spatial reasoning on the CSS2 visual box model. In *Proceedings of the 21st National Conference on Artificial Intelligence (AAAI'06)*, pages 1313–1318, Boston, Massachusetts, July 2006.
- [25] W. Gatterbauer, P. Bohunsky, M. Herzog, B. Krüpl, and B. Pollak. Towards domain-independent information extraction from web tables. In *Proceedings of the 16th International World Wide Web Conference (WWW'07)*, pages 71–80, Banff, Canada, May 2007.
- [26] S. Handschuh, S. Staab, and F. Ciravegna. S-CREAM — Semi-automatic CREATION of Metadata. In *Proceedings of the 13th International Conference on Knowledge Engineering and Knowledge Management (EKAW'02)*, pages 358–372, Sigüenza, Spain, October 2002.
- [27] W. Holzinger, B. Krüpl, and M. Herzoge. Using ontologies for extracting product features from web pages. In *Proceedings of the Fifth International Semantic Web Conference (ISWC'06)*, pages 286–299, Athens, Georgia, November 2006.
- [28] P.G. Ipeirotis, L. Gravano, and M. Sahami. Probe, count, and classify: categorizing hidden web databases. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data (SIGMOD'01)*, pages 67–78, Santa Barbara, California, May 2001.
- [29] P. Kogut and W. Holmes. AeroDAML: Applying information extraction to generate DAML annotations from web pages. In *Proceedings of the of the First International Conference on Knowledge Capture (K-CAP'01) Workshop on Knowledge Markup and Semantic Annotation*, Victoria, British Columbia, 2001.
- [30] K. Lerman, L. Getoor, S. Minton, and C. Knoblock. Using the structure of web sites for automatic segmentation of tables. In *Proceedings of the 2004 ACM SIGMOD International Conference on Management of Data (SIGMOD'04)*, pages 119–130, Paris, France, 2004.

- [31] S. Lim and Y. Ng. An automated approach for retrieving hierarchical data from HTML tables. In *Proceedings of the Eighth International Conference on Information and Knowledge Management (CIKM'99)*, pages 466–474, Kansas City, Missouri, November 1999.
- [32] R. Navigli, P. Velardi, A. Cucchiarelli, and F. Neri. Quantitative and qualitative evaluation of the OntoLearn ontology learning system. In *Proceedings of the 20th International Conference on Computational Linguistics*, pages 1043–1050, Geneva, Switzerland, August 2004.
- [33] N.F. Noy, M. Sintek, S. Decker, M. Crubezy, R.W. Fergerson, and M. Musen. Creating semantic web contents with Protégè-2000. *IEEE Intelligent Systems*, 16(2):60–71, March/April 2001.
- [34] A. Pivk. Automatic ontology generation from web tabular structures. *AI Communications*, 19(1):83–85, 2006.
- [35] A. Pivk, P. Cimiano, and Y. Sure. From tables to frames. In *Proceedings of the Third International Semantic Web Conference (ISWC'04)*, pages 166–181, Hiroshima, Japan, November 2004.
- [36] B. Popov, A. Kiryakov, D. Ognyanoff, D. Manov, and A. Kirilov. KIM — a semantic platform for information extraction and retrieval. *Natural Language Engineering*, 10(3-4):375–392, 2004.
- [37] P. Spyns, D. Oberle, R. Volz, J. Zheng, M. Jarrar, Y. Sure, R. Studer, and R. Meersman. OntoWeb—a semantic web community portal. In *Proceedings of the 5th International Conference on Practical Aspects of Knowledge Management (PAKM'02)*, pages 189–200, Vienna, Austria, December 2002.
- [38] K.-C. Tai. The tree-to-tree correction problem. *Journal of the ACM*, 26(3):422–433, 1979.
- [39] C. Tao and D.W. Embley. Automatic hidden-web table interpretation by sibling page comparison. In *Proceedings of the 26th International Conference on Conceptual Modeling (ER'07)*, pages 560–581, Auckland, New Zealand, November 2007.
- [40] C. Tao and D.W. Embley. Seed-based generation of personalized bio-ontologies for information extraction. In *Proceedings of the First International Workshop on Conceptual Modelling for Life Sciences Applications (CMLSA'07)*, pages 74–84, Auckland, New Zealand, November 2007.
- [41] A. Tengli, Y. Yang, and N.L. Ma. Learning table extraction from examples. In *Proceedings the 20th International Conference on Computational Linguistics (COLING'04)*, pages 987–993, Geneva, Switzerland, August 2004.

- [42] C.A. Thompson, M.E. Califf, and R.J. Mooney. Active learning for natural language parsing and information extraction. In *Proceedings of 16th International Conference on Machine Learning*, pages 406–414, Bled, Slovenia, June 1999.
- [43] Y.A. Tijerino, D.W. Embley, D.W. Lonsdale, Y. Ding, and G. Nagy. Toward ontology generation from tables. *World Wide Web: Internet and Web Information Systems*, 8(3):251–285, September 2004.
- [44] M. Vargas-Vera, E. Motta, J. Domingue, M. Lanzoni, A. Stutt, and F. Ciravegna. MnM: Ontology driven semi-automatic and automatic support for semantic markup. In *Proceedings of the 13th International Conference on Knowledge Engineering and Knowledge Management (EKAW'02)*, pages 379–391, Siguenza, Spain, October 2002.
- [45] X. Wang. *Tabular Abstraction, Editing, and Formatting*. PhD thesis, Univeristy of Waterloo, 1996.
- [46] Y. Wang and J. Hu. A machine learning based approach for table detection on the web. In *Proceedings of the 11th International Conference on World Wide Web (WWW'02)*, pages 242–250, Honolulu, Hawaii, May 2002.
- [47] Y. Wang, J. Völker, and P. Haase. Towards semi-automatic ontology building supported by large-scale knowledge acquisition. In *AAAI Fall Symposium On Semantic Web for Collaborative Knowledge Acquisition*, volume FS-06-06, pages 70–77, Arlington, Virginia, October 2006.
- [48] W. Yang. Identifying syntactic differences between two programs. *Software Practice and Experience*, 21(7):739–755, 1991.
- [49] R. Zanibbi, D. Blostein, and J.R. Cordy. A survey of table recognition. *International Journal of Document Analysis and Recognition*, 7(1):1–16, 2004.
- [50] Y. Zhai and B. Liu. Web data extraction based on partial tree alignment. In *Proceedings of the 14th International Conference on World Wide Web (WWW'05)*, pages 76–85, Chiba, Japan, May 2005.